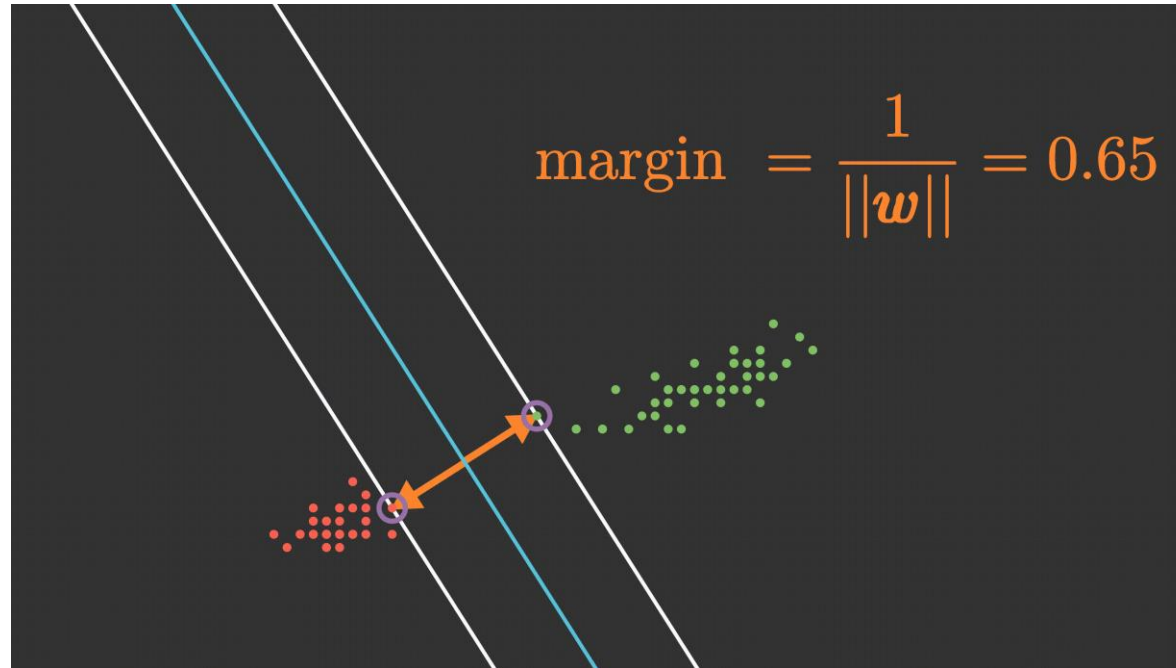


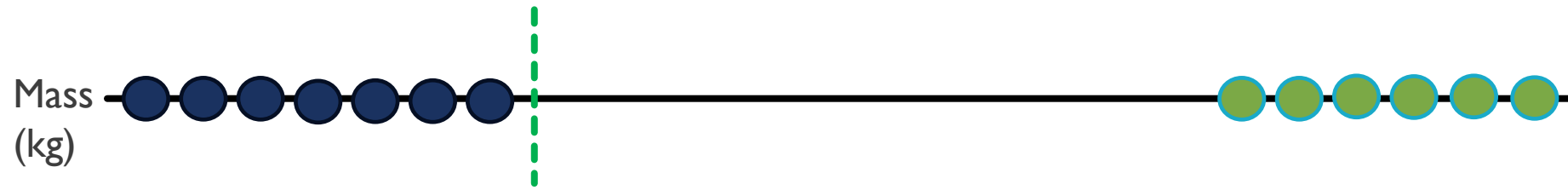
CT-562 MACHINE LEARNING

NED University of Engineering & Technology

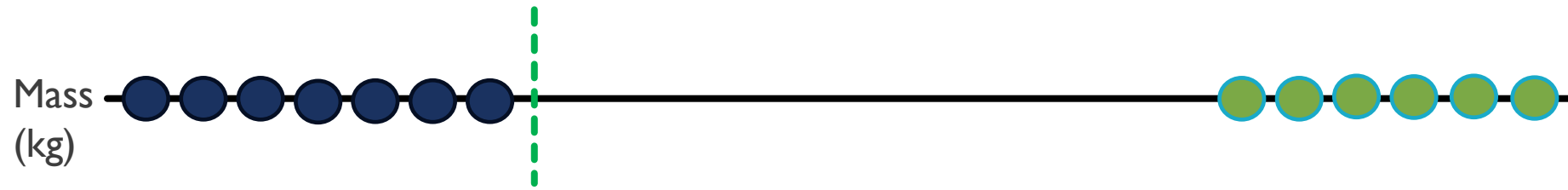


SUPPORT VECTOR MACHINE

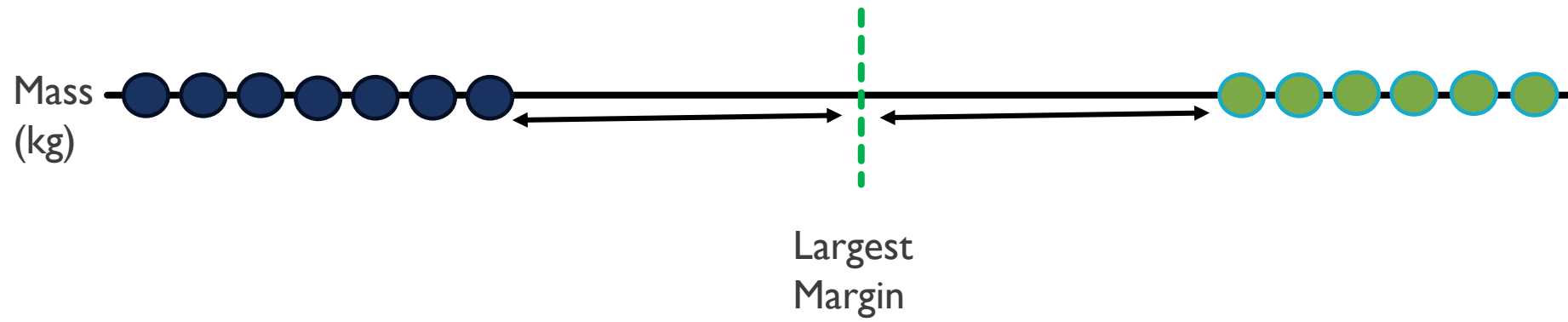
DATA



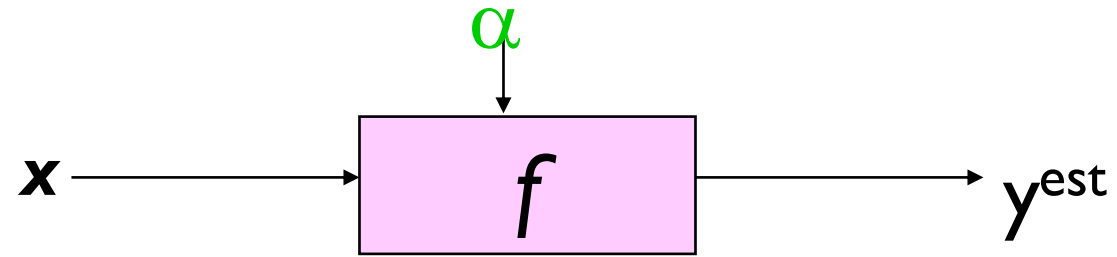
DATA



DATA

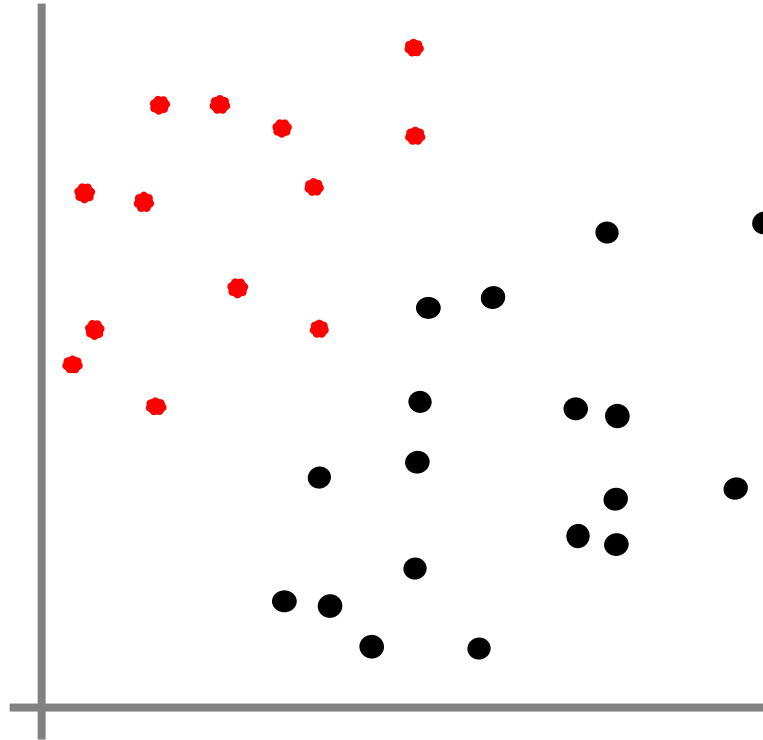


DATA



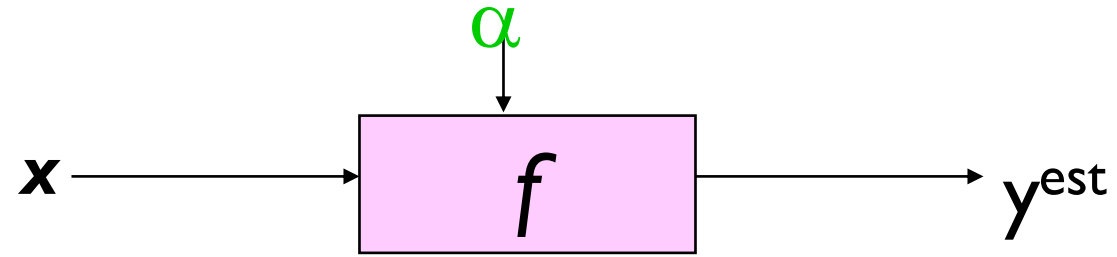
$$f(x, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



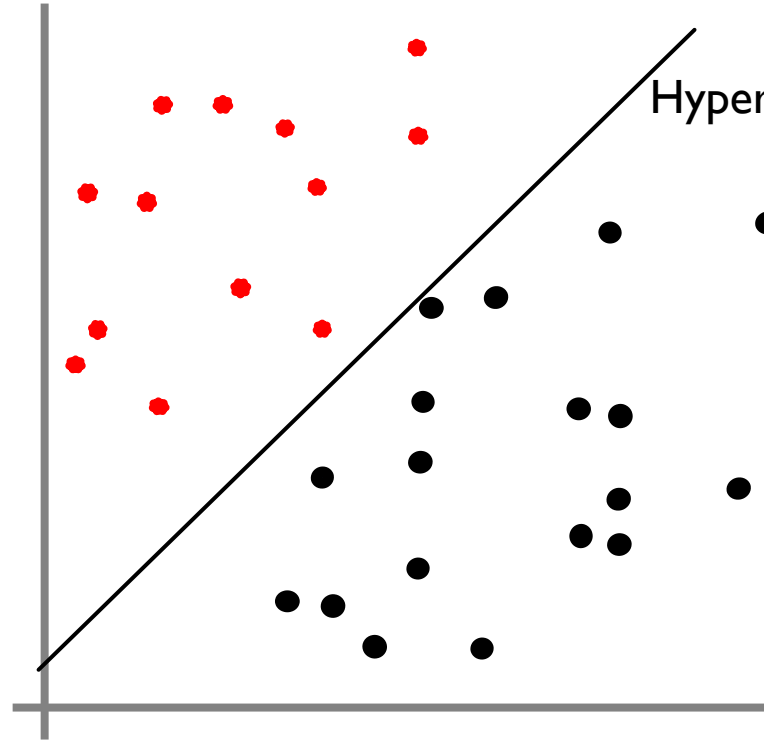
How would you
classify this data?

HYPERPLANE



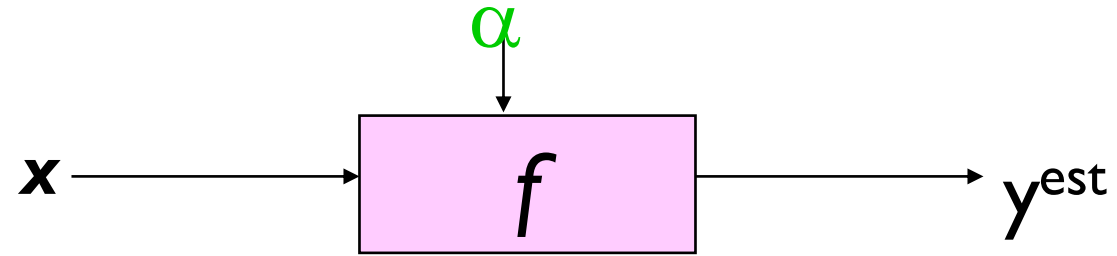
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



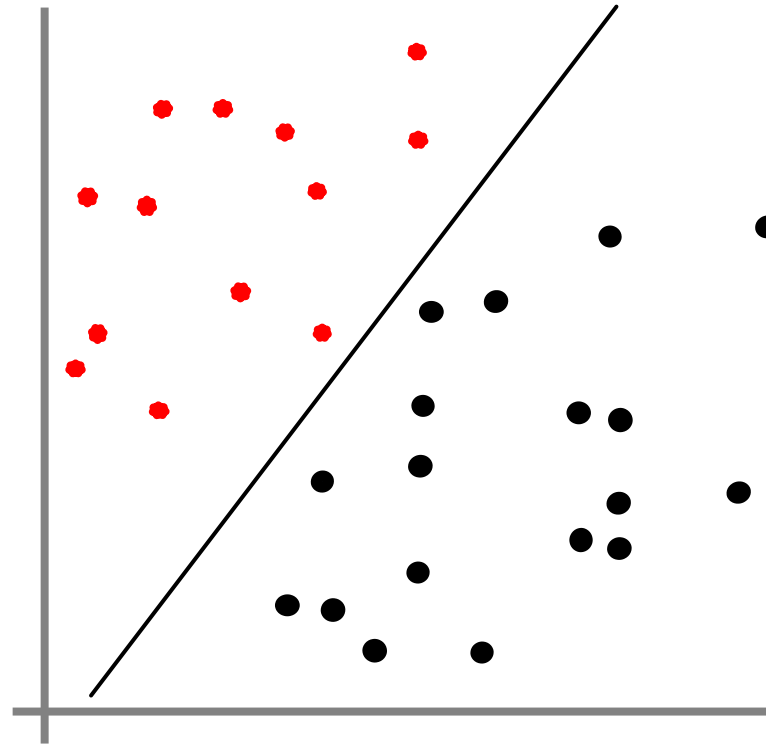
How would you
classify this data?

HYPERPLANE



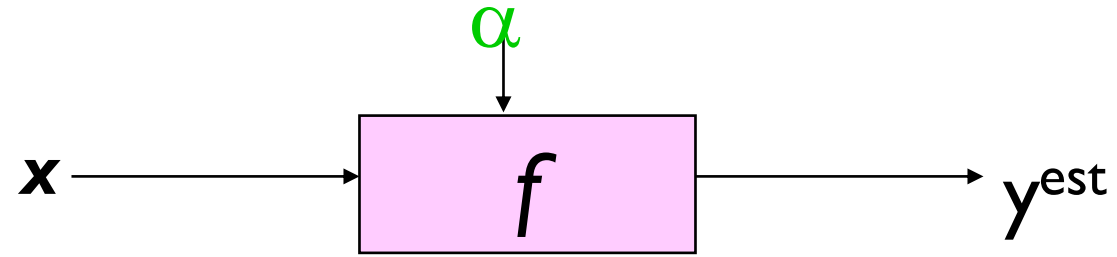
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



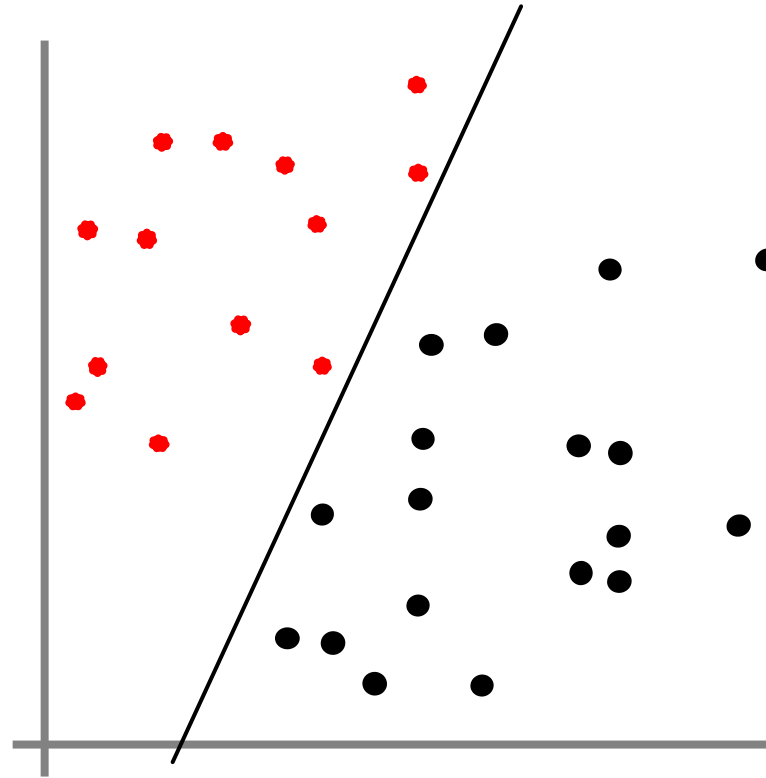
How would you
classify this data?

HYPERPLANE



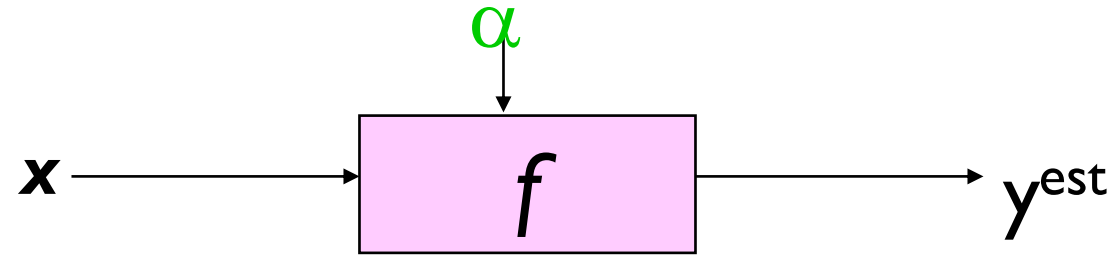
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

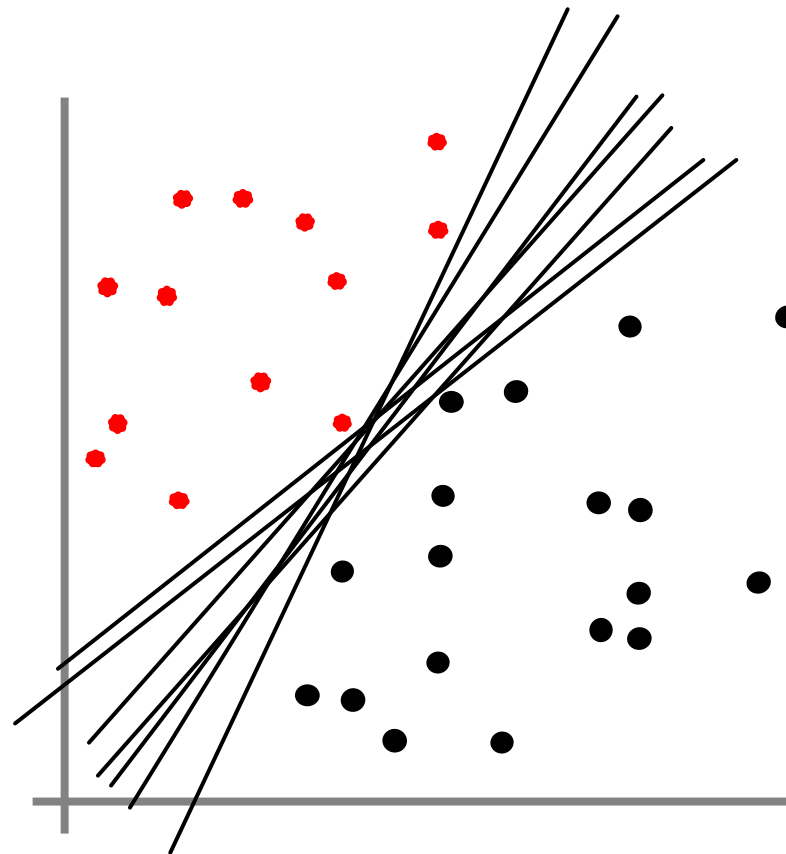


How would you
classify this data?

HYPERPLANE



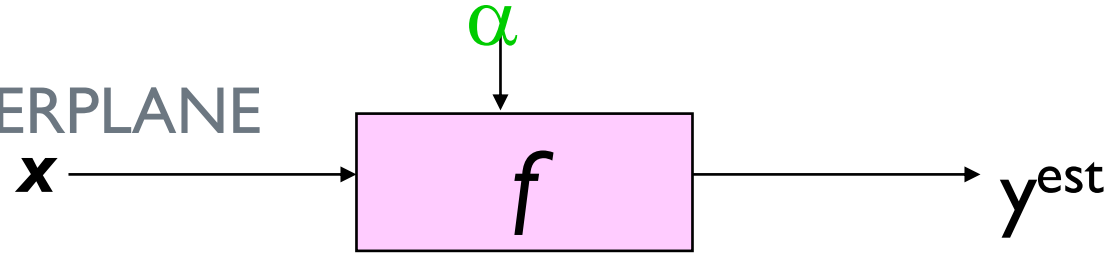
- denotes +1
- denotes -1



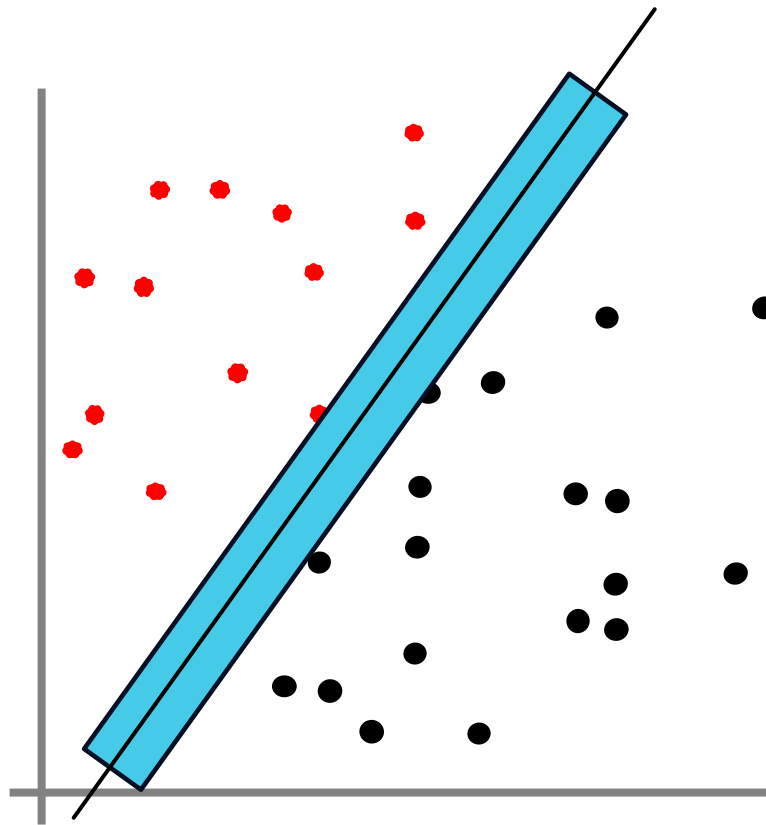
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

How would you
classify this data?

MAXIMAL MARGIN HYPERPLANE



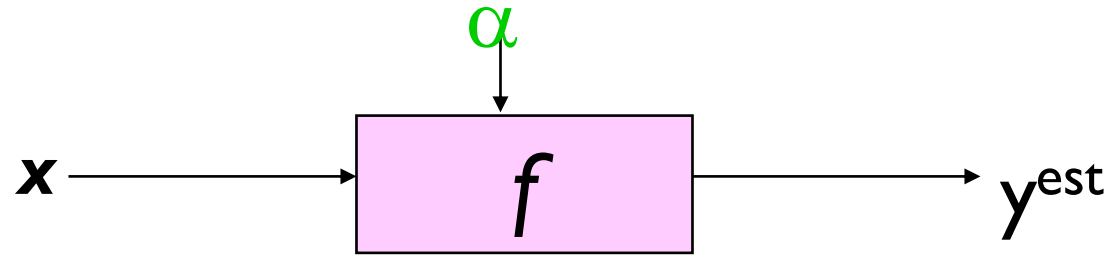
- denotes +1
- denotes -1



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

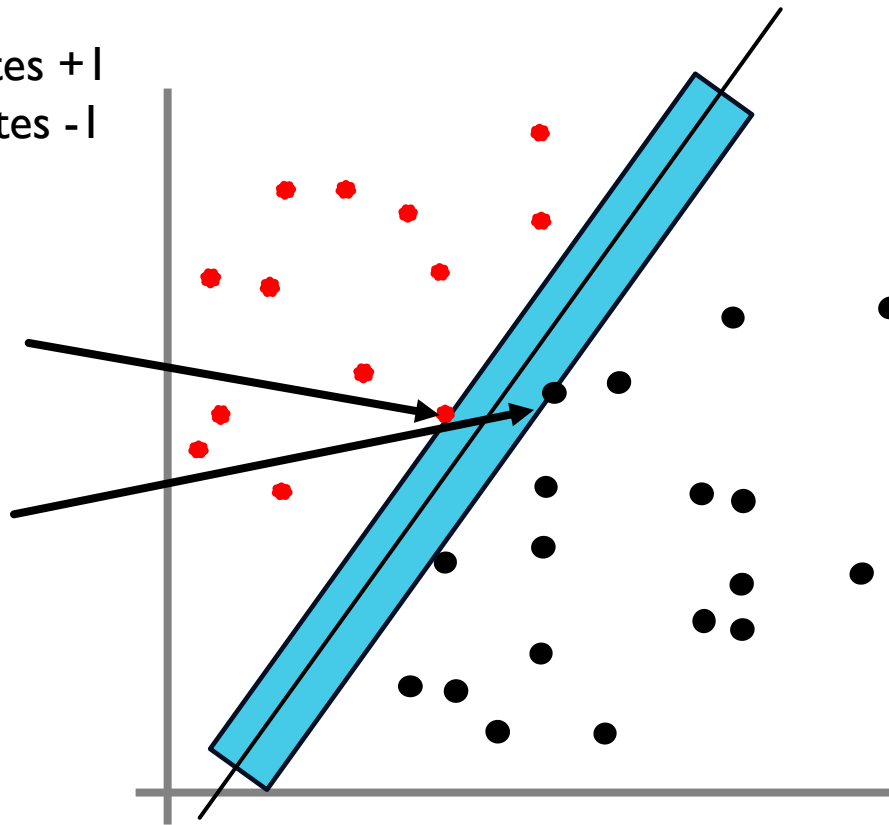
LINEAR CLASSIFIERS



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

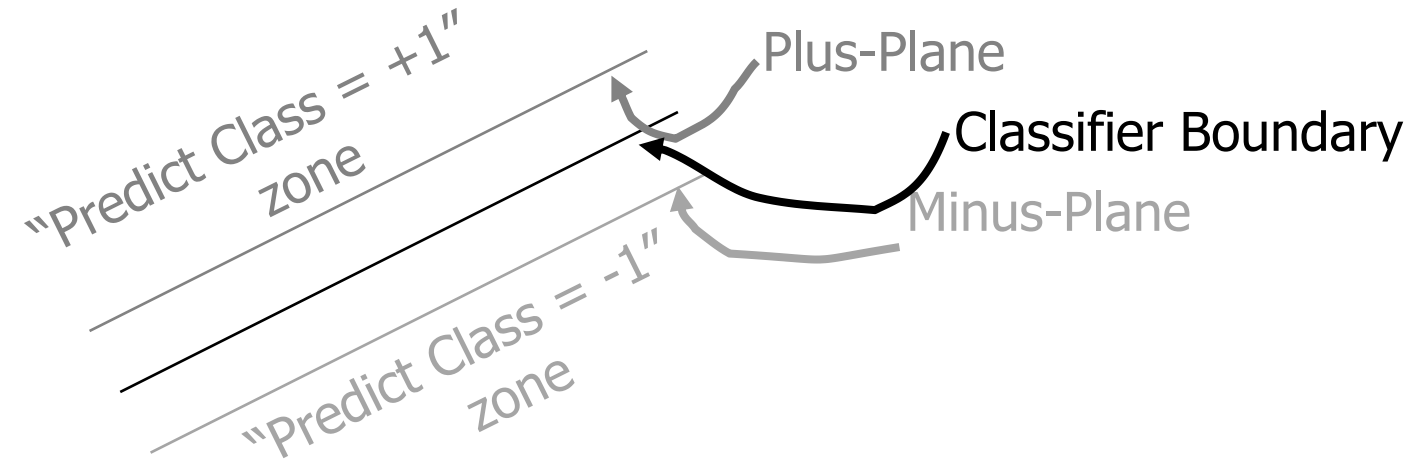
Support Vectors are those datapoints that the margin pushes up against



The maximum margin linear classifier is the linear classifier with the, maximum margin. This is the simplest kind of SVM (Called an LSVM)

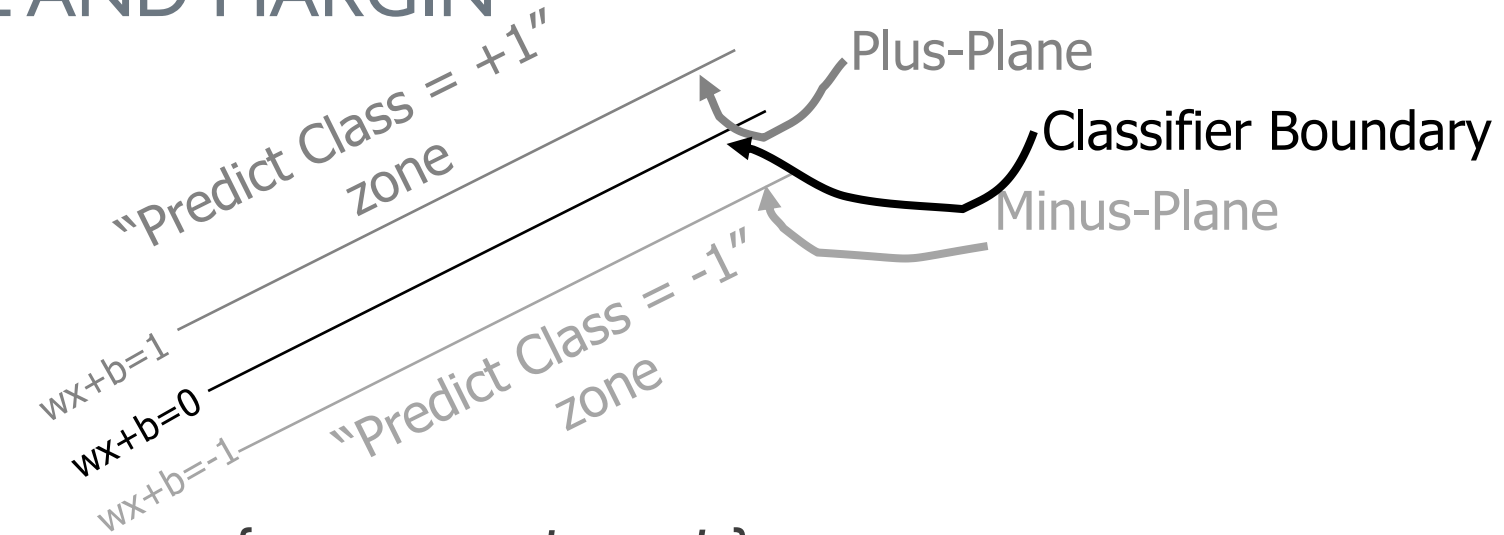
Linear SVM

SPECIFYING A LINE AND MARGIN



- How do we represent this mathematically?
- ...in m input dimensions?

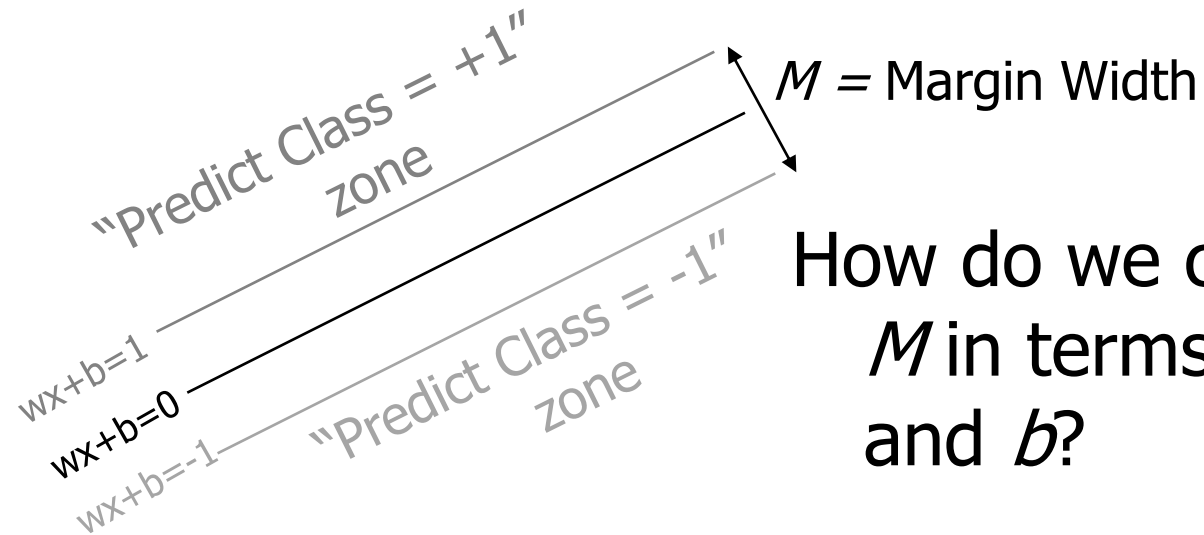
SPECIFYING A LINE AND MARGIN



- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

Classify as..	+1	if	$\mathbf{w} \cdot \mathbf{x} + b \geq 1$
	-1	if	$\mathbf{w} \cdot \mathbf{x} + b \leq -1$
	Universe explodes	if	$-1 < \mathbf{w} \cdot \mathbf{x} + b < 1$

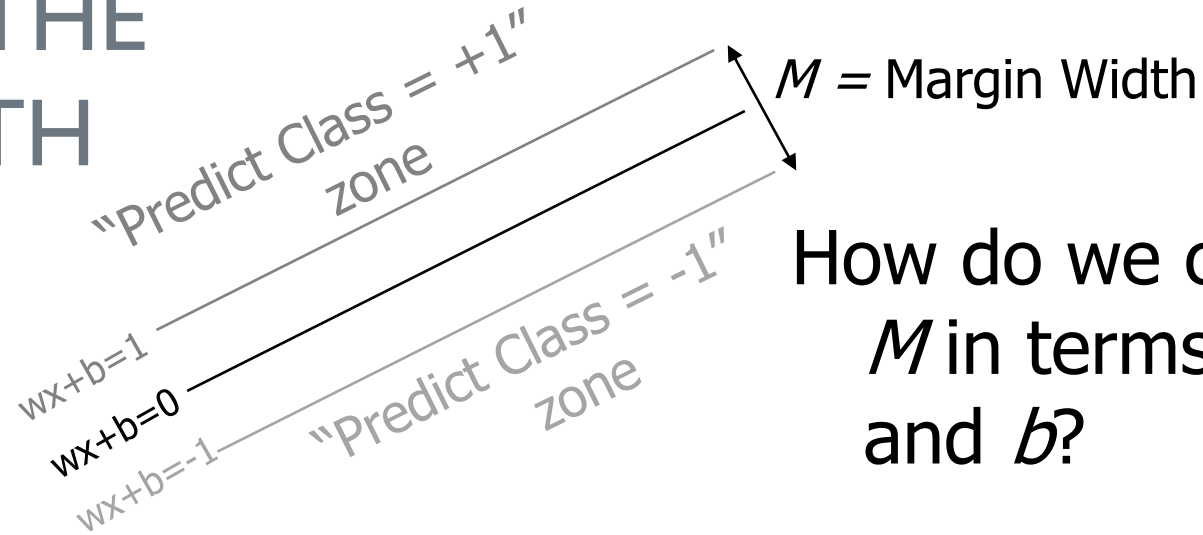
COMPUTING THE MARGIN WIDTH



How do we compute M in terms of \mathbf{w} and b ?

- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

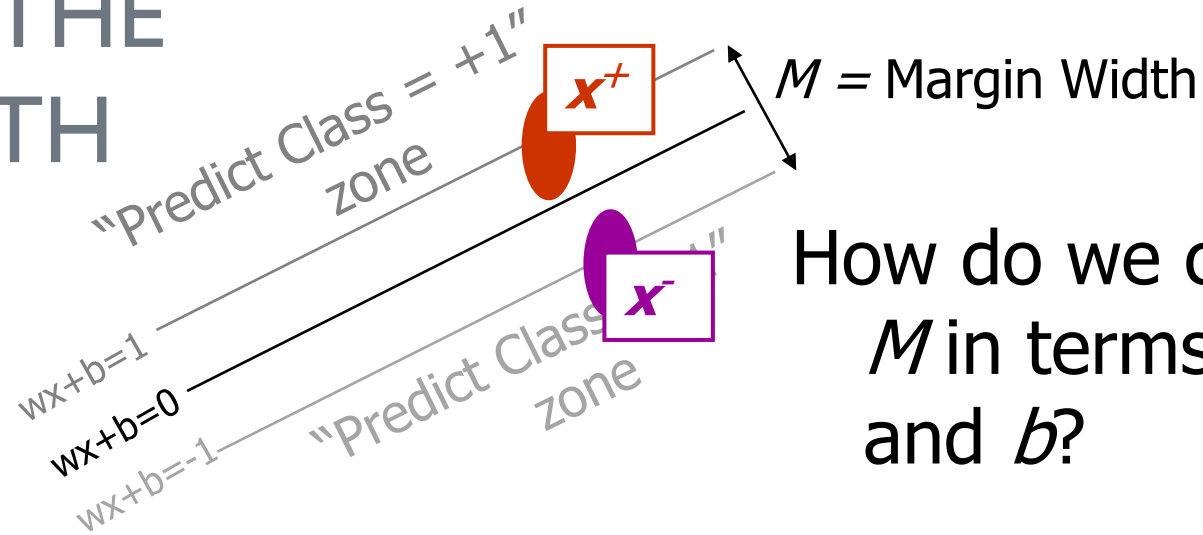
COMPUTING THE MARGIN WIDTH



How do we compute M in terms of \mathbf{w} and b ?

- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

COMPUTING THE MARGIN WIDTH

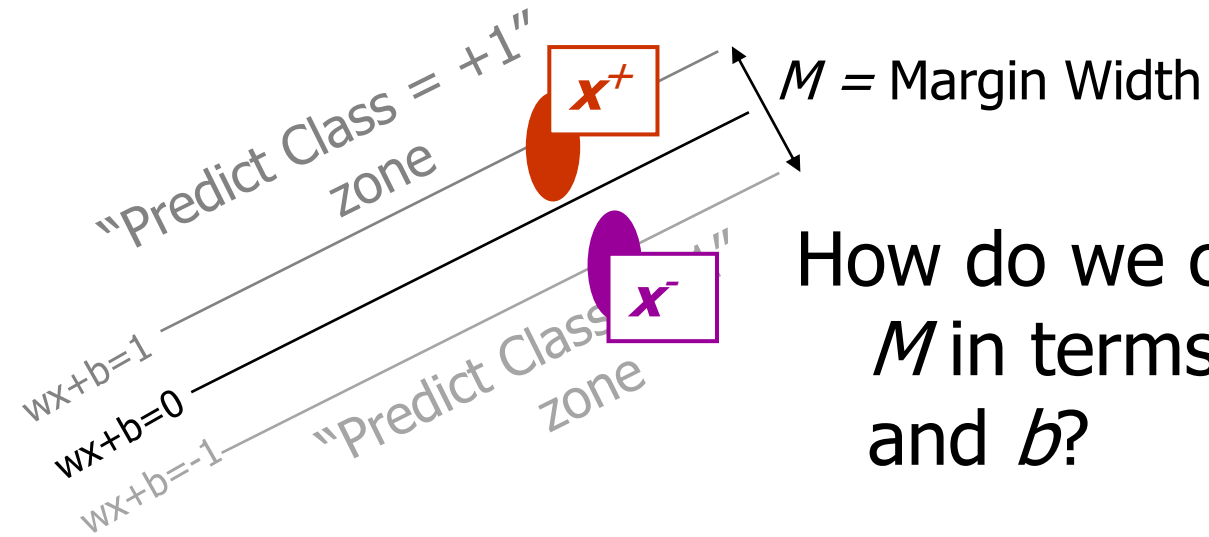


How do we compute M in terms of \mathbf{w} and b ?

- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector \mathbf{w} is perpendicular to the Plus Plane
- Let \mathbf{x}^- be any point on the minus plane
- Let \mathbf{x}^+ be the closest plus-plane-point to \mathbf{x}^- .

Any location in \mathbb{R}^m : not necessarily a datapoint

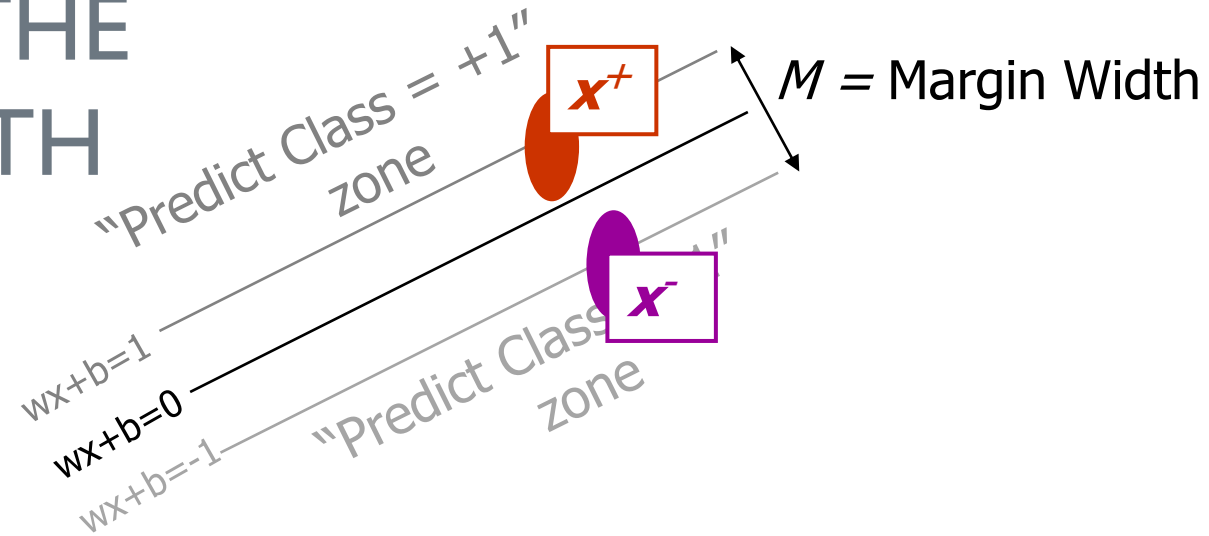
COMPUTING THE MARGIN WIDTH



How do we compute M in terms of \mathbf{w} and b ?

- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector \mathbf{w} is perpendicular to the Plus Plane
- Let \mathbf{x}^- be any point on the minus plane
- Let \mathbf{x}^+ be the closest plus-plane-point to \mathbf{x}^- .

COMPUTING THE MARGIN WIDTH

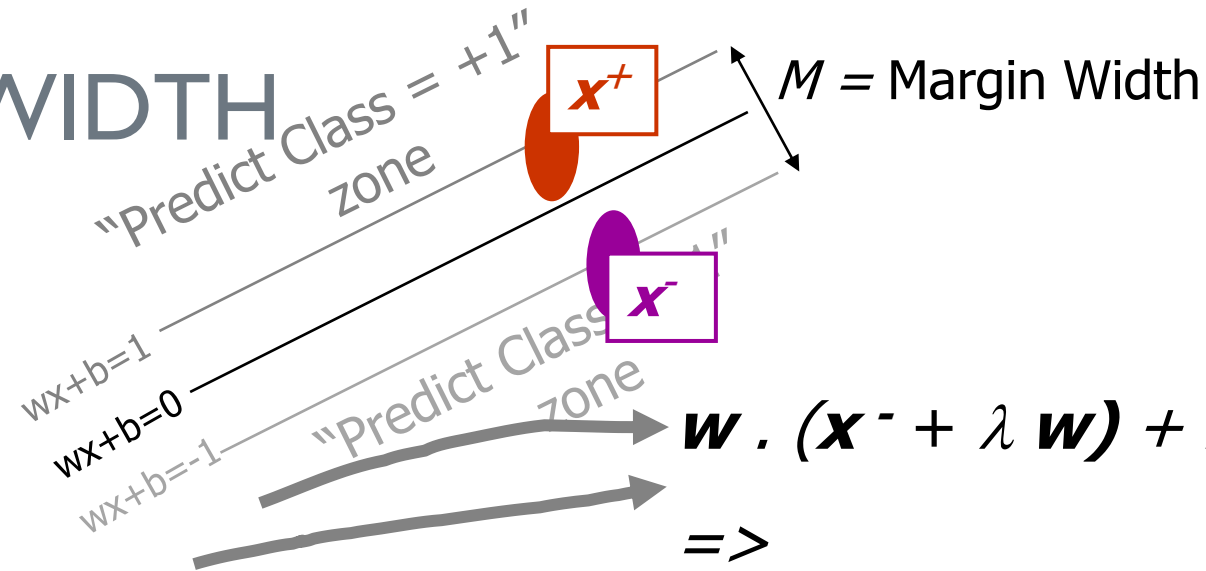


What we know:

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

It's now easy to get M in terms of w and b

COMPUTING THE MARGIN WIDTH



What we know:

- $w . x^+ + b = +1$
- $w . x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

It's now easy to get M in terms of w and b

$$w . (x^- + \lambda w) + b = 1$$

\Rightarrow

$$w . x^- + b + \lambda w . w = 1$$

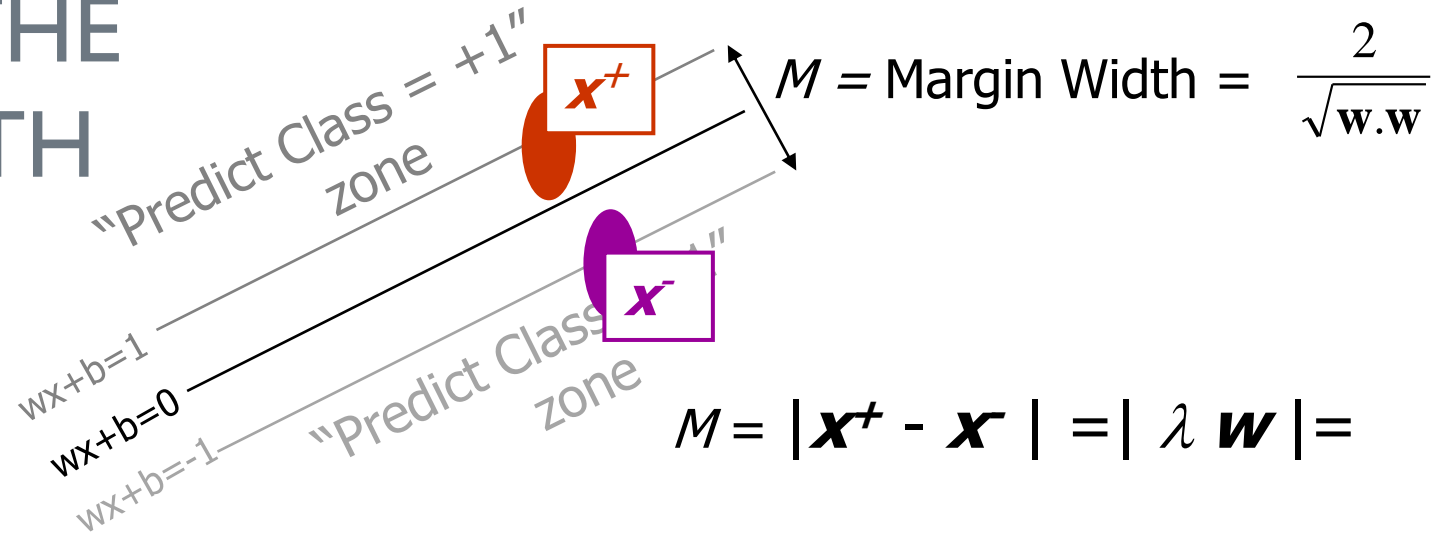
\Rightarrow

$$-1 + \lambda w . w = 1$$

\Rightarrow

$$\lambda = \frac{2}{w.w}$$

COMPUTING THE MARGIN WIDTH



What we know:

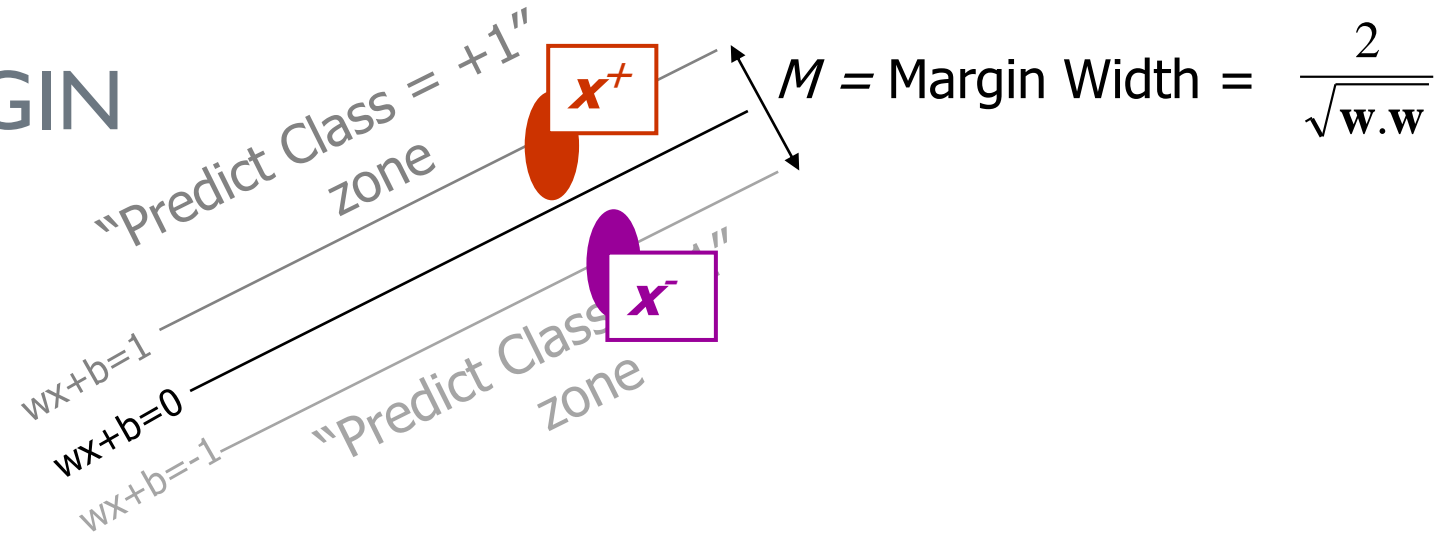
- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$
- $\lambda = \frac{2}{\sqrt{w \cdot w}}$

$$M = |x^+ - x^-| = |\lambda w| =$$

$$= \lambda |w| = \lambda \sqrt{w \cdot w}$$

$$= \frac{2\sqrt{w \cdot w}}{w \cdot w} = \frac{2}{\sqrt{w \cdot w}}$$

LEARNING THE MAXIMUM MARGIN CLASSIFIER



Given a guess of \mathbf{w} and b we can

- Compute whether all data points in the correct half-planes
- Compute the width of the margin

So now we just need to write a program to search the space of \mathbf{w} 's and b 's to find the widest margin that matches all the datapoints.

THE MATH

- Training instances
 - $\mathbf{x} \in \mathcal{R}^n$
 - $y \in \{-1, 1\}$
- Decision function
 - $f(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$
 - $\mathbf{w} \in \mathcal{R}^n$
 - $b \in \mathcal{R}$
- Find \mathbf{w} and b that
 - Perfectly classify training instances
 - Assuming linear separability
 - Maximize margin

THE MATH

- For perfect classification, we want
 - $y_i (\langle w, x_i \rangle + b) \geq 0$ for all i
 - Why?
- To maximize the margin, we want
 - w that minimizes $|w|^2$

STRENGTHS OF SVMs

- Good generalization in theory
- Good generalization in practice
- Work well with few training instances
- Efficient algorithms

WHAT IF SURFACE IS NON-LINEAR?

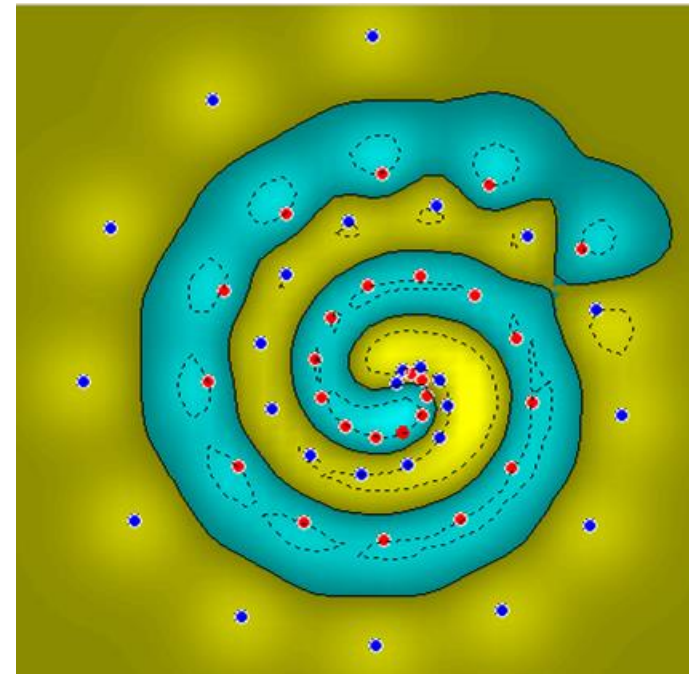
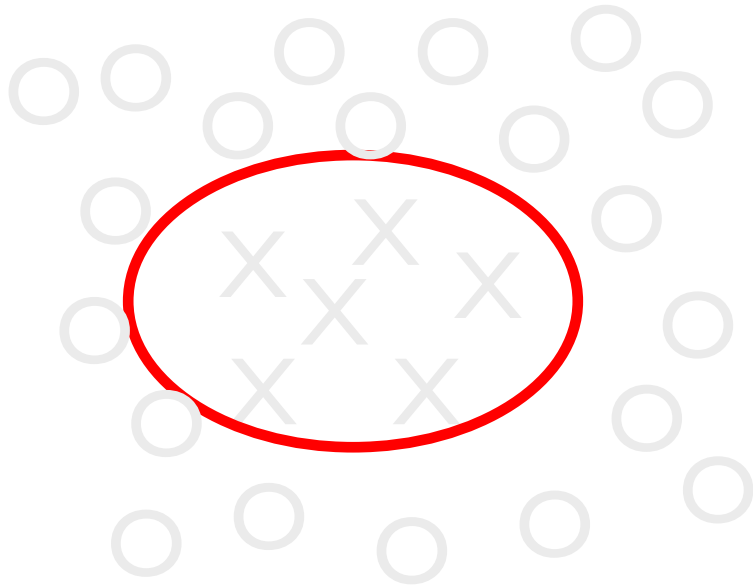


Image from <http://www.atrandomresearch.com/iclass/>

WHAT IF SURFACE IS NON-LINEAR?

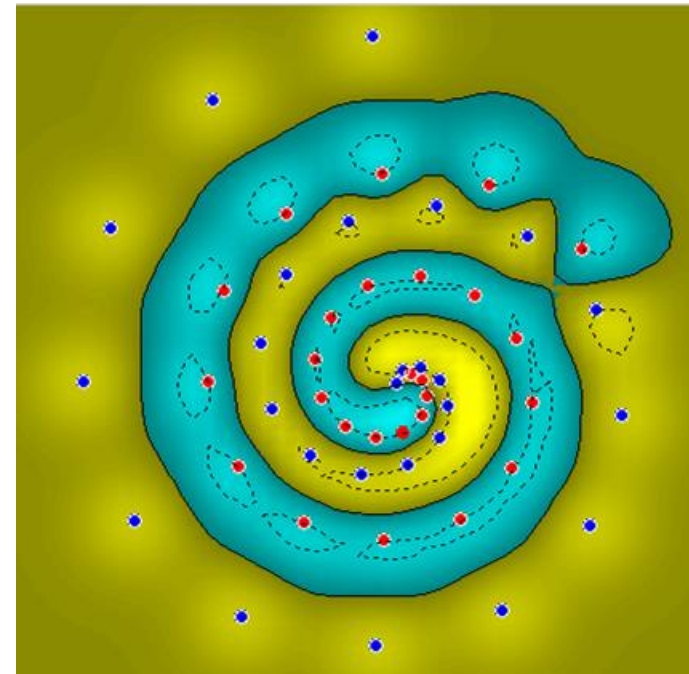
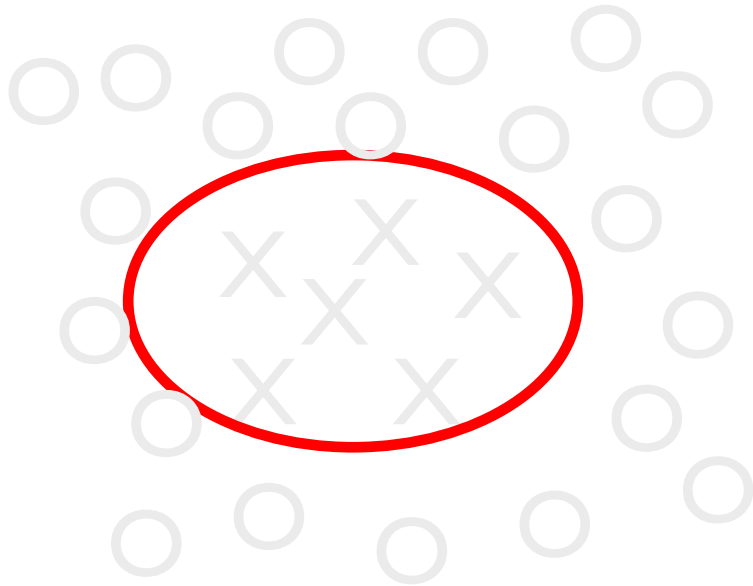
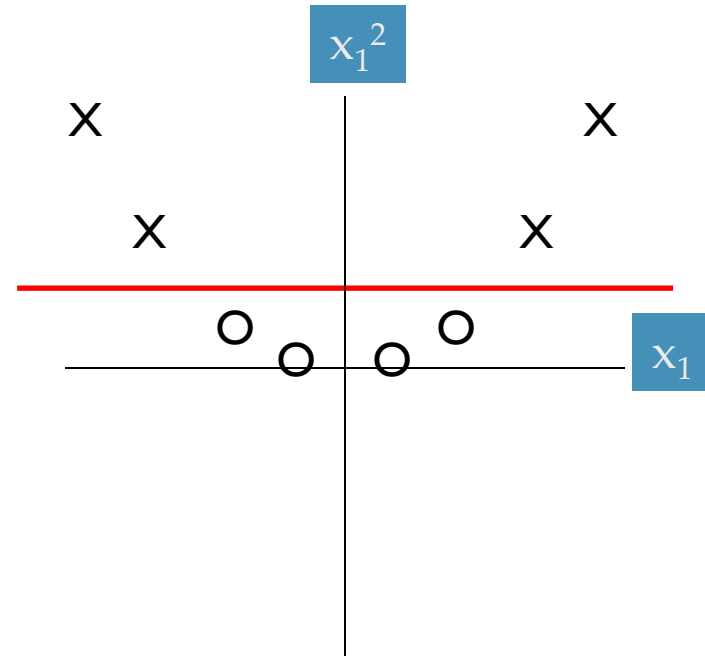
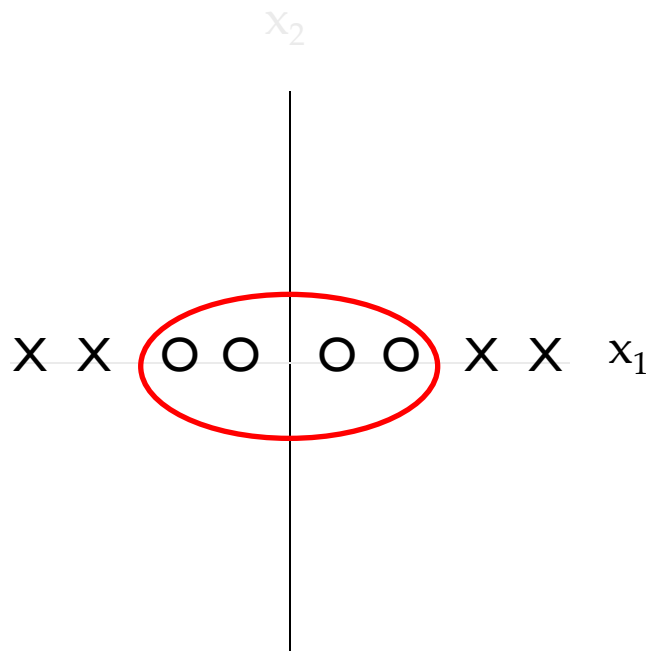
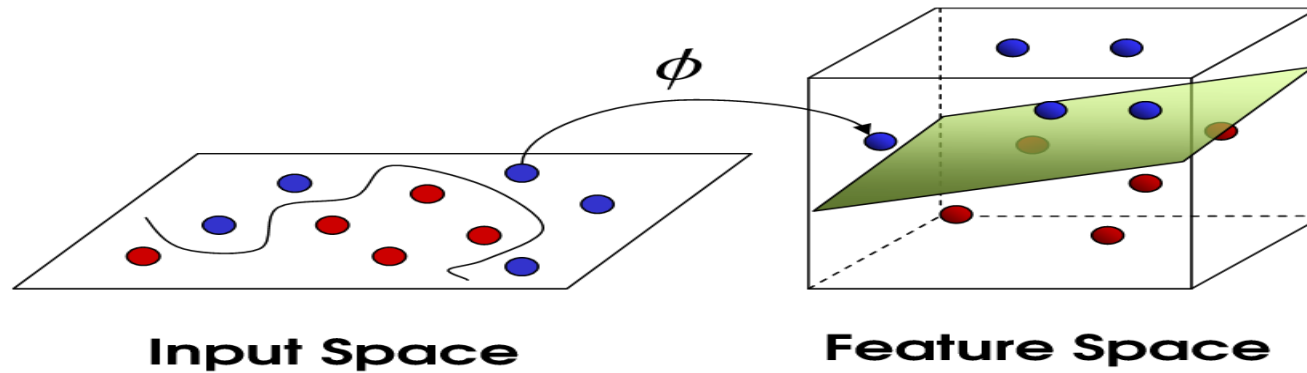


Image from <http://www.atrandomresearch.com/iclass/>

WHEN LINEAR SEPARATORS FAIL



MAPPING INTO A NEW FEATURE SPACE



$$\Phi : x \rightarrow X = \Phi(x)$$

$$\Phi(x_1, x_2) = (x_1, x_2, x_1^2, x_2^2, x_1x_2)$$

- Rather than run SVM on x_i , run it on $\Phi(x_i)$
- Find non-linear separator in input space
- What if $\Phi(x_i)$ is really big?
- Use kernels to compute it implicitly!

KERNELS

- Find kernel K such that
 - $K(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle$
- Computing $K(x_1, x_2)$ should be efficient, much more so than computing $\Phi(x_1)$ and $\Phi(x_2)$
- Use $K(x_1, x_2)$ in SVM algorithm rather than $\langle x_1, x_2 \rangle$
- Remarkably, this is possible

THE POLYNOMIAL KERNEL

- $K(x_1, x_2) = \langle x_1, x_2 \rangle^2$
 - $x_1 = (x_{11}, x_{12})$
 - $x_2 = (x_{21}, x_{22})$
- $\langle x_1, x_2 \rangle = (x_{11}x_{21} + x_{12}x_{22})$
- $\langle x_1, x_2 \rangle^2 = (x_{11}^2x_{21}^2 + x_{12}^2x_{22}^2 + 2x_{11}x_{12}x_{21}x_{22})$
- $\Phi(x_1) = (x_{11}^2, x_{12}^2, \sqrt{2}x_{11}x_{12})$
- $\Phi(x_2) = (x_{21}^2, x_{22}^2, \sqrt{2}x_{21}x_{22})$
- $K(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle$

THE POLYNOMIAL KERNEL

- $\Phi(x)$ contains all monomials of degree d
- Useful in visual pattern recognition
- Number of monomials
 - 16x16 pixel image
 - 10^{10} monomials of degree 5
- Never explicitly compute $\Phi(x)$!
- Variation - $K(x_1, x_2) = (\langle x_1, x_2 \rangle + 1)^2$

A FEW GOOD KERNELS

- Dot product kernel

- $K(x_1, x_2) = \langle x_1, x_2 \rangle$

- Polynomial kernel

- $K(x_1, x_2) = \langle x_1, x_2 \rangle^d$ (Monomials of degree d)

- $K(x_1, x_2) = (\langle x_1, x_2 \rangle + 1)^d$ (All monomials of degree $1, 2, \dots, d$)

- Gaussian kernel

- $K(x_1, x_2) = \exp(-\|x_1 - x_2\|^2 / 2\sigma^2)$

- Radial basis functions

- Sigmoid kernel

- $K(x_1, x_2) = \tanh(\langle x_1, x_2 \rangle + \theta)$

- Neural networks

- Establishing “kernel-hood” from first principles is non-trivial

THE KERNEL TRICK

“Given an algorithm which is formulated in terms of a positive definite kernel K_1 , one can construct an alternative algorithm by replacing K_1 with another positive definite kernel K_2 ”

➤ SVMs can use the kernel trick

EXOTIC KERNELS

- Strings
- Trees
- Graphs
- The hard part is establishing kernel-hood



CONCLUSION

- SVMs find optimal linear separator
- The kernel trick makes SVMs non-linear learning algorithms



THANK YOU